

Introduction to probability

Reference: Bertsekas Chapter I

Outline

- Probability and conditional probability
- Probability models
- Set-theoretic concepts and notation
- Kolmogorov's rules and the multiplication rule
- Probability calculus and basic counting rule
- Probabilities from data
- Rule of total probability
- Independence
- Fallacy of mistaken compounding
- Bayes Rule

Probability basics

If A denotes some event, then $P(A)$ is the probability that this event occurs:

- $P(\text{coin lands heads}) = 0.5$
- $P(\text{rainy day in Ireland}) = 0.85$
- $P(\text{cold day in Hell}) = 0.0000001$

And so on.

Probability basics

Some probabilities are estimated from direct experience over the long run:

- $P(\text{newborn baby is a boy}) = \frac{106}{206}$
- $P(\text{death due to car accident}) = \frac{11}{100,000}$
- $P(\text{death due to any cause}) = 1$

Probability basics

Some probabilities are estimated from direct experience over the long run:

- $P(\text{newborn baby is a boy}) = \frac{106}{206}$
- $P(\text{death due to car accident}) = \frac{11}{100,000}$
- $P(\text{death due to any cause}) = 1$

Others are synthesized from our best judgments about unique events:

- $P(\text{Apple stock goes up after next earnings call}) = 0.54$
- $P(\text{Djokovic wins next US Open}) = 0.4$ (6 to 4 odds)
- etc.

Conditional probability

A conditional probability is the chance that one thing happens, given that some other thing has already happened.

A great example is a weather forecast: if you look outside this morning and see gathering clouds, you might assume that rain is likely and carry an umbrella.

We express this judgment as a conditional probability: e.g. “the conditional probability of rain this afternoon, given clouds this morning, is 60%.”

Conditional probability

In stats we write this a bit more compactly:

- $P(\text{rain this afternoon} \mid \text{clouds this morning}) = 0.6$
- That vertical bar means “given” or “conditional upon.”
- The thing on the left of the bar is the event we’re interested in.
- The thing on the right of the bar is our knowledge, also called the “conditioning event” or “conditioning variable”: what we believe or assume to be true.

$P(A \mid B)$: “the probability of A, given that B occurs.”

Conditional probability

Conditional probabilities are how we express judgments in a way that reflects our partial knowledge.

- You just gave *Sherlock* a high rating. What's the conditional probability that you will like *The Imitation Game* or *Tinker Tailor Soldier Spy*?
- You just bought organic dog food on Amazon. What's the conditional probability that you will also buy a GPS-enabled dog collar?
- You follow Lionel Messi (@leomessi) on Instagram. What's the conditional probability that you will respond to a suggestion to follow Cristiano Ronaldo (@cristiano) or Gareth Bale (@garethbale11)?

Conditional probability

A really important fact is that conditional probabilities are *not symmetric*:

$$P(A \mid B) \neq P(B \mid A)$$

As a quick counter-example, let the events A and B be as follows:

- A: “you can dribble a basketball”
- B: “you play in the NBA”

Conditional probability

- A: “you can dribble a basketball”
- B: “you play in the NBA”



Clearly $P(A \mid B) = 1$: every NBA player can dribble a basketball.

Conditional probability

- A: “you can dribble a basketball”
- B: “you play in the NBA”



But $P(B | A)$ is nearly zero!

Uncertain outcomes and probability models

An uncertain outcome (more formally called a “random process”) has two key properties:

1. The set of possible outcomes, called the *sample space*, is known beforehand.
2. The particular outcome that occurs is *not* known beforehand.

We denote the sample space as Ω , and some particular element of the sample space as $\omega \in \Omega$.

Uncertain outcomes and probability models

Examples:

1. NBA finals, Golden State vs. Toronto:

$$\Omega = \{4-0, 4-1, 4-2, 4-3, 3-4, 2-4, 1-4, 0-4\}$$

2. Temperature in degrees F in Austin on a random day:

$$\Omega = [10, 115]$$

3. Number of no-shows on an AA flight from Austin to DFW:

$$\Omega = \{0, 1, 2, \dots, N_{\text{seats}}\}$$

4. Poker hand

$$\Omega = \text{all possible five-card deals from a 52-card deck}$$

Uncertain outcomes and probability models

An event is a *subset of the sample space*, i.e. $A \subset \Omega$. For example:

1. NBA finals, Golden State vs. Toronto. Let A be the event “Toronto wins”. Then

$$A = \{3-4, 2-4, 1-4, 0-4\} \subset \Omega$$

2. Austin weather. Let A be the event “cooler than 90 degrees”. Then

$$A = [10, 90) \subset [10, 115]$$

3. Flight no-shows. Let A be “more than 5 no shows”:

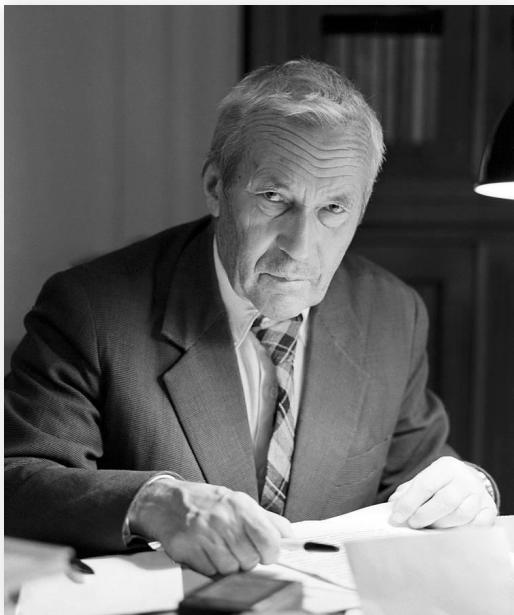
$$A = \{6, 7, 8, \dots, N_{\text{seats}}\}$$

Some set-notation reminders

We need some basic set-theory concepts to make sense of probability, since the sample space Ω is a set, and since “events” are subsets of Ω .

- Union: $A \cup B = \{\omega : \omega \in A \text{ or } \omega \in B\}$
- Intersection: $A \cap B = \{\omega : \omega \in A \text{ and } \omega \in B\}$
- Complement: $A^C = \tilde{A} = \{\omega : \omega \notin A\}$
- Difference/relative complement: $A \setminus B = \{\omega : \omega \in A, \omega \notin B\}$
- Disjointness: A and B are disjoint if $A \cap B = \emptyset$ (the empty set).

Kolmogorov's axioms (baby version)



“Obey my rules,
filthy capitalists.”

Consider an uncertain outcome with sample space Ω . “Probability” $P(\cdot)$ is a set function that maps Ω to the real numbers, such that:

1. Non-negativity: For any event $A \subset \Omega$, $P(A) \geq 0$.
2. Normalization: $P(\Omega) = 1$ and $P(\emptyset) = 0$.
3. Finite additivity: If A and B are disjoint, then $P(A \cup B) = P(A) + P(B)$.

Not that intuitive! Notice no mention of frequencies...

Some optional technical points

Point I: the “non-baby” version of Kolmogorov's third axiom is actually something called *countable additivity* (versus “finite additivity”).

- Consider a countable sequence of events A_1, A_2, A_3, \dots , each $A_i \subset \Omega$.
- Suppose the events are all disjoint: $A_i \cap A_j = \emptyset$ for $i \neq j$.
- Countable additivity says that

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

where

$$\bigcup_{i=1}^{\infty} A_i = A_1 \cup A_2 \cup A_3 \cup \dots$$

Some optional technical points

Point 2: our definition of probability assumes that $P(A)$ is well defined for any $A \subset \Omega$. If Ω is a finite or countable set (e.g. the integers), then $P(A)$ is always well defined.

But if Ω is uncountably infinite (e.g. the real numbers), then $P(A)$ is not necessarily well defined for all possible subsets $A \subset \Omega$.

- Wild but true! It is possible to define bizarre sets for which there is *no meaningful notion of that set's size*. (If you care: Google “non-measurable set” and “Banach-Tarski paradox”).
- But it's hard to construct such crazy sets. Every “normal” set you might care about (intervals, unions of intervals) has a well-defined size. Technically speaking, these are called the *Borel sets*.

Quick summary of terms

- Uncertain outcome/“random process”: we know the possibilities ahead of time, just not the specific one that occurs.
- Sample space: the set of possible outcomes.
- Event: a subset of the sample space.
- Probability: a function that maps events to real numbers and that obeys Kolmogorov's axioms.

OK, so how do we actually *calculate* probabilities?

The discrete uniform distribution and the counting rule

Suppose our sample space Ω is a finite set consisting of N elements $\omega_1, \dots, \omega_N$.

Suppose further that $P(\omega_i) = 1/N$: each outcome is equally likely, i.e. we have a *discrete uniform distribution* over possible outcomes.

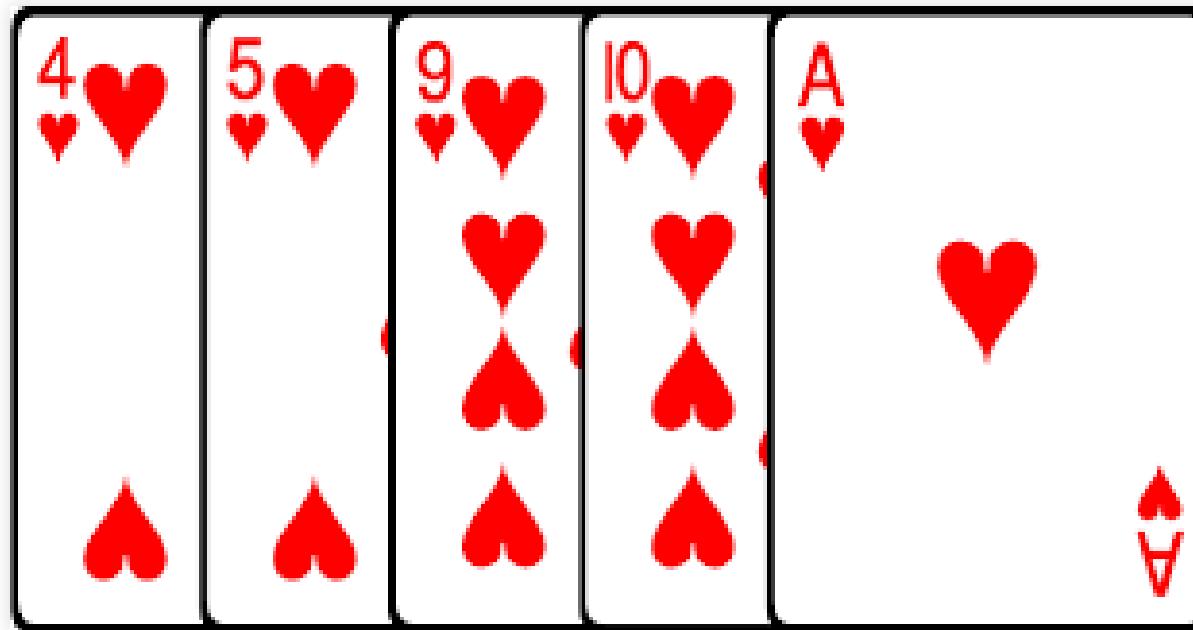
Then for each set $A \subset \Omega$,

$$P(A) = \frac{|A|}{N} = \frac{\text{Number of elements in } A}{\text{Number of elements in } \Omega}$$

That is, to compute $P(A)$, we just need to count how many elements are in A .

The counting rule: example

Someone deals you a five-card poker hand from a 52-card deck.
What is the probability of a flush (all five cards the same suit)?



Note: this is a very historically accurate illustration of probability, given its origins among bored French aristocrats!

The counting rule: example

Let's use the counting rule.

- Our sample space has $N = \binom{52}{5} = 2,598,960$ possible poker hands, each one equally likely.
- How many possible flushes are there? Let's start with hearts:
 - There are 13 hearts.
 - To make a flush with hearts, you need any 5 of these 13 cards.
 - Thus there are $\binom{13}{5} = 1287$ possible flushes with hearts.
- The same argument works for all four suits, so there are $4 \times 1287 = 5,148$ flushes. Thus

$$P(\text{flush}) = \frac{|A|}{|\Omega|} = \frac{5148}{2598960} = 0.00198079$$

Probability calculus

The “probability calculus” provides a set of rules for calculating probabilities. *These aren't axioms:* they can be derived from Kolmogorov's axioms.

1. $P(A^C) = 1 - P(A)$

(Why? Because $A \cup A^C = \Omega$, and $P(\Omega) = 1$.)

2. If $A \subset B$, then $P(A) \leq P(B)$.

(Why? Write B as $B = A \cup (B \setminus A)$ and use finite additivity.)

3. $P(B \setminus A) = P(B) - P(A \cap B)$.

(Why? **Your turn!**)

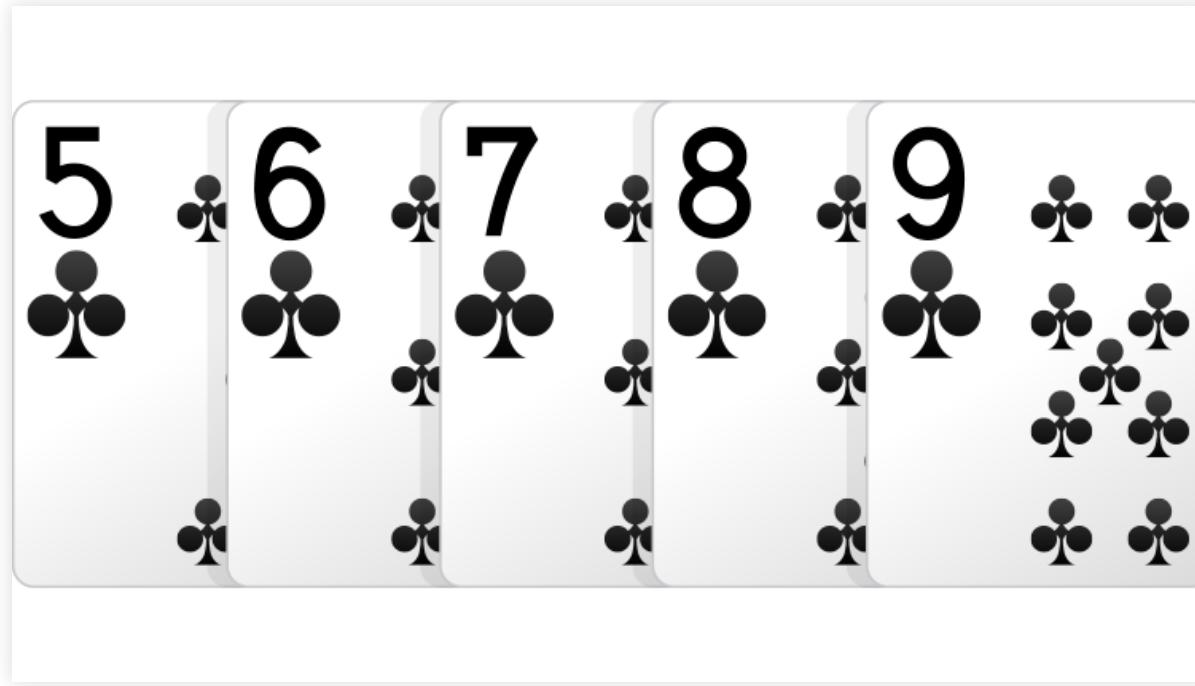
4. Addition rule: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

(Why? **Your turn!**)

Quick example

Again, someone deals you a five-card poker hand. What is the probability of either a straight (five cards in a row, e.g. 3-4-5-6-7) or a flush (all cards the same suit)?

Note: these aren't mutually exclusive, since you might draw a hand that is both a straight AND a flush (e.g. 5-6-7-8-9 of clubs).



Quick example

If all 2,598,960 poker hands are equally likely, then using the counting rule:

- $P(\text{flush}) = 0.00198079$ (5,148 possible flushes)
- $P(\text{straight}) = 0.00392465$ (10,200 possible straights.)
- $P(\text{straight and flush}) = 0.0000153908$ (40 possible straight flushes)

So by the addition rule:

$$\begin{aligned}P(\text{straight or flush}) &= P(\text{straight}) + P(\text{flush}) - P(\text{straight AND flush}) \\&= 0.00392465 + 0.00198079 - 0.0000153908 \\&= 0.005890049\end{aligned}$$

The multiplication rule

We've met Kolmogorov's three axioms, together with several rules we can derive from these axioms.

There's one final axiom for conditional probability, often called the *multiplication rule*. Let $P(A, B) = P(A \cap B)$ be the *joint probability* that both A and B happen. Then:

$$P(A \mid B) = \frac{P(A, B)}{P(B)} .$$

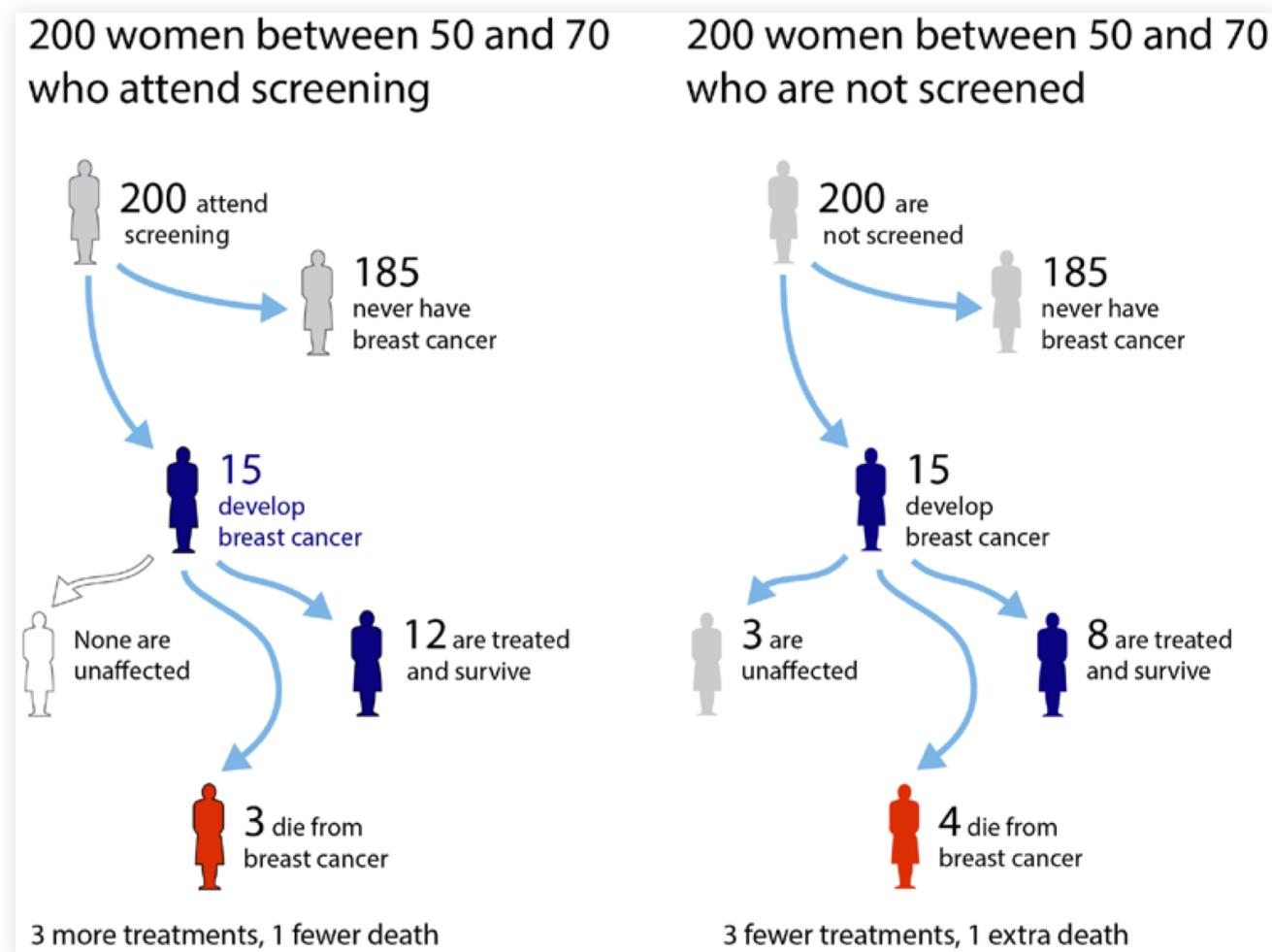
Or equivalently:

$$P(A, B) = P(A \mid B) \cdot P(B) .$$

This is an axiom: it cannot be proven from Kolmogorov's rules.

Conditional probability

Let's see why this axiom makes sense. (Figure courtesy David Speigelhalter and Jenny Gage.)



Conditional probability

Suppose a woman goes for regular screening (left branch). What is $P(\text{survive} \mid \text{cancer})$?

- Among screened women, 15 will get cancer, on average.
- Of those 15, 12 are treated and survive, on average.
- Thus intuitively, we should have $P(\text{survive} \mid \text{cancer}) = 12/15 = 0.8$

We get the same answer using the rule for conditional probabilities:

$$\begin{aligned} P(S|C) &= \frac{P(S, C)}{P(C)} = \frac{12/200}{15/200} \\ &= \frac{12}{15} \end{aligned}$$

Conditional probabilities from data

Consider the story of Abraham Wald and the “missing”“ WWII bombers.



Conditional probabilities from data

- B-17 bombers returning from their missions in WWII often had damage: on the fuselage, across the wings, on the engine block, and sometimes even near the cockpit.
- At some point, a clever data-minded person had the idea of analyzing the distribution of these hits over the surface of the returning planes.
- The thinking was that, if you could find patterns in where the B-17s were taking enemy fire, you could figure out where to reinforce them with extra armor, to improve survivability.
- You couldn't reinforce them everywhere, or they would be too heavy to fly.

Conditional probabilities from data

Suppose we saw data on returning bombers that looked like this:

| Location | Number of planes |
|-----------------------|------------------|
| Engine | 53 |
| Cockpit area | 65 |
| Fuel system | 96 |
| Wings, fuselage, etc. | 434 |

Naive answer: of 648 returning planes, 434 (68%) were hit on the fuselage.

$$P(\text{hit on wings or fuselage} \mid \text{returns home}) = 0.68$$

So bulk up the fuselage!

Conditional probabilities from data

But that's the right answer to the wrong question! We need the inverse probability

$$P(\text{returned home} \mid \text{hit on wings or fuselage})$$

Wald used some fancy math to reconstruct the *joint frequency distribution* over damage type and mission result:

| | Returned | Shot down |
|-----------------------|----------|-----------|
| Engine | 53 | 57 |
| Cockpit area | 65 | 46 |
| Fuel system | 96 | 16 |
| Wings, fuselage, etc. | 434 | 33 |

Conditional probabilities from data

This gives us:

$$P(\text{returns safely} \mid \text{hit on wings or fuselage}) = \frac{434}{434 + 33} \approx 0.93.$$

Pretty high.

Conditional probabilities from data

On the other hand, of the 110 planes that had taken damage to the engine, only 53 only returned safely. Therefore

$$P(\text{returns safely} \mid \text{hit on engine}) = \frac{53}{53 + 57} \approx 0.48.$$

Similarly,

$$P(\text{returns safely} \mid \text{hit on cockpit area}) = \frac{65}{65 + 46} \approx 0.59.$$

Moral of the story: make sure you're focusing on the right conditional probability. And remember that $P(A \mid B) \neq P(B \mid A)$!

Conditional probabilities from data

The same math that Abraham Wald used to analyze bullet holes on B-17s also underpins the modern digital economy of films, television, music, and social media.

Netflix, Hulu, and other video-streaming services all use this same math to examine what shows their users are watching, and apply the results of their number-crunching to recommend new shows.

Many companies do the same:

- Amazon for products
- New York Times for new stories
- Google for web pages
- etc

Conditional probabilities from data

Suppose that you're designing the movie-recommendation algorithm for Netflix, and you have access to the entire Netflix database, showing which customers have liked which films.

Your goal is to leverage this vast data resource to make automated, personalized movie recommendations.

You decide to start with an easy case: assessing how probable it is that a user will like the film *Saving Private Ryan* (event *A*), given that the same user has liked the HBO series *Band of Brothers* (event *B*).

This is almost certainly a good bet: both are epic war dramas about the Normandy invasion and its aftermath. Therefore, you might think: job done! Recommend away.

Conditional probabilities from data

But keep in mind that you want to be able to do this kind of thing automatically.

Key insight: frame the problem in terms of conditional probability.

Suppose we learn that Linda liked film B , but hasn't yet seen film A .

- What if $P(\text{likes } A \mid \text{likes } B) = 0.8$? Then A is a good recommendation! Based on Linda liking B , we know there's an 80% chance she'll like A .
- But if $P(\text{likes } A \mid \text{likes } B) = 0.02$, she probably won't like A , given our knowledge of how she reacted to B .

Conditional probabilities hold the key to understanding individualized preferences. So how can we learn $P(\text{likes } A \mid \text{likes } B)$?

Conditional probabilities from data

Solution: go to the data! Suppose your database on 5 million subscribers like Linda reveals the following pattern:

| | Liked <i>Band of Brothers</i> | Didn't like |
|----------------------------------|-------------------------------|-------------|
| Liked <i>Saving Private Ryan</i> | 2.8 million | 0.3 million |
| Didn't like | 0.7 million | 1.2 million |

Then

$$P(\text{liked Saving Private Ryan} \mid \text{liked Band of Brothers}) = \frac{2.8 \text{ million}}{3.5 \text{ million}} = 0.8.$$

Result: a good recommendation with no human in the loop.

Conditional probabilities from data

Let's try an example ourselves in `predimed_intro.R` from the class website.

We'll use data from a large experiment to estimate:

- $P(\text{cardiac event} \mid \text{Mediterranean diet})$
- $P(\text{cardiac event} \mid \text{control diet})$

The rule of total probability

Consider the following data on complication rates at a maternity hospital in Cambridge, England:

| | Easier deliveries | Harder deliveries | Overall |
|----------------|--------------------------|--------------------------|----------------|
| Senior doctors | 0.052 | 0.127 | 0.076 |
| Junior doctors | 0.067 | 0.155 | 0.072 |

Would you rather have a junior or senior doctor?

The rule of total probability

Consider the following data on complication rates at a maternity hospital in Cambridge, England:

| | Easier deliveries | Harder deliveries | Overall |
|----------------|--------------------------|--------------------------|----------------|
| Senior doctors | 0.052 | 0.127 | 0.076 |
| Junior doctors | 0.067 | 0.155 | 0.072 |

Would you rather have a junior or senior doctor?

Simpson's paradox. Senior doctors have:

- lower complication rates for easy cases.
- lower complication rates for hard cases.
- higher complication rates overall! (7.6% versus 7.2%). **Why?**

The rule of total probability

Let's see the table with number of deliveries performed (in parentheses):

| | Easier deliveries | Harder deliveries | Overall |
|----------------|--------------------------|--------------------------|----------------|
| Senior doctors | 0.052 (213) | 0.127 (102) | 0.076 (315) |
| Junior doctors | 0.067 (3169) | 0.155 (206) | 0.072 (3375) |

The rule of total probability

Let's see the table with number of deliveries performed (in parentheses):

| | Easier deliveries | Harder deliveries | Overall |
|----------------|--------------------------|--------------------------|----------------|
| Senior doctors | 0.052 (213) | 0.127 (102) | 0.076 (315) |
| Junior doctors | 0.067 (3169) | 0.155 (206) | 0.072 (3375) |

Now we see what's going on:

- Most of the deliveries performed by junior doctors are easier cases, where complication rates are lower overall.
- The senior doctors, meanwhile, work a much higher fraction of the harder cases.

The rule of total probability

It turns out the math of Simpson's paradox can be understood a lot more deeply in terms of something called the *rule of total probability*, or the mixture rule.

This rule sounds fancy, but is actually quite simple.

It says to divide and conquer: **the probability of any event is the sum of the probabilities for all the different ways in which the event can happen.** Really just Kolmogorov's third rule in disguise!

The rule of total probability

Let's see this rule in action for the hospital data.

There are two types of deliveries: easy and hard. So:

$$P(\text{complication}) = P(\text{easy and complication}) + P(\text{hard and complication}) .$$

Now use the rule for conditional probabilities to each joint probability on the right-hand side:

$$\begin{aligned} P(\text{complication}) &= P(\text{easy}) \cdot P(\text{complication} \mid \text{easy}) \\ &\quad + P(\text{hard}) \cdot P(\text{complication} \mid \text{hard}) . \end{aligned}$$

The rule of total probability says that overall probability is a weighted average—a **mixture**—of the two conditional probabilities

The rule of total probability

For senior doctors we get

$$P(\text{complication}) = \frac{213}{315} \cdot 0.052 + \frac{102}{315} \cdot 0.127 = 0.076.$$

And for junior doctors, we get

$$P(\text{complication}) = \frac{3169}{3375} \cdot 0.067 + \frac{206}{3375} \cdot 0.155 = 0.072.$$

This is a lower *marginal* or *overall* probability of a complication, even though junior doctors have higher *conditional* probabilities of a complication in all scenarios.

Synonyms: overall probability = total probability = marginal probability

The rule of total probability

You can see why these are called *marginal* probabilities if you go back to the Abraham Wald example:

| | Returned | Shot down |
|-----------------------|-----------------|------------------|
| Engine | 53 | 57 |
| Cockpit area | 65 | 46 |
| Fuel system | 96 | 16 |
| Wings, fuselage, etc. | 434 | 33 |

The rule of total probability

You can see why these are called *marginal* probabilities if you go back to the Abraham Wald example:

| | Returned | Shot down |
|-----------------------|-----------------|------------------|
| Engine | 0.066 | 0.071 |
| Cockpit area | 0.081 | 0.058 |
| Fuel system | 0.120 | 0.020 |
| Wings, fuselage, etc. | 0.542 | 0.042 |

Divide by the total number (800) to turn these into joint probabilities...

The rule of total probability

You can see why these are called *marginal* probabilities if you go back to the Abraham Wald example:

| | Returned | Shot down | Marginal |
|-----------------------|-----------------|------------------|-----------------|
| Engine | 0.066 | 0.071 | 0.137 |
| Cockpit area | 0.081 | 0.058 | 0.139 |
| Fuel system | 0.120 | 0.020 | 0.140 |
| Wings, fuselage, etc. | 0.542 | 0.042 | 0.584 |
| Marginal | 0.809 | 0.191 | 1 |

Now calculate the overall probabilities for each individual type of event, and put those in the margins of the table.

The rule of total probability

Here's the formal statement of the rule. Let Ω be any sample space, and let B_1, B_2, \dots, B_N be a *partition* of Ω —that is, a set of events such that:

$$P(B_i, B_j) = 0 \text{ for any } i \neq j, \quad \text{and} \quad \sum_{i=1}^N P(B_i) = 1.$$

Now consider any event A . Then

$$P(A) = \sum_{i=1}^N P(A, B_i) = \sum_{i=1}^N P(B_i) \cdot P(A \mid B_i).$$

Example: drug surveys

Virginia Delaney-Black and her colleagues at Wayne State University gave an anonymous survey to teenagers in Detroit:

- 432 teens were asked whether they had used various drugs.
- Of these 432 teens, 211 agreed to give a hair sample.
- Therefore, for these 211 respondents, the researchers could compare people's answers with an actual drug test.
- Hair samples were analyzed in the aggregate: no hair sample could be traced back to an individual survey or teen.

Citation: V. Delaney-Black et. al. "Just Say I Don't: Lack of Concordance Between Teen Report and Biological Measures of Drug Use." *Pediatrics* 165:5, pp. 887-93 (2010)

Example: drug surveys

The two sets of results were strikingly different.

- Of the 211 teens who provided a hair sample, only a tiny fraction of them (0.7%) admitted to having used cocaine.
- But when the hair samples were analyzed in the lab, 69 of them (33.7%) came back positive for cocaine use.

Example: drug surveys

The two sets of results were strikingly different.

- Of the 211 teens who provided a hair sample, only a tiny fraction of them (0.7%) admitted to having used cocaine.
- But when the hair samples were analyzed in the lab, 69 of them (33.7%) came back positive for cocaine use.

And the parents lied, too:

- The researchers also asked the parents whether they had used cocaine themselves.
- Only 6.1% said yes.
- But 28.3% of the parents' hair samples came back positive.

Example: drug surveys

Remember:

- these people were guaranteed anonymity
- they wouldn't be arrested or fired for saying admitting drug use
- they willingly agreed to provide a hair sample that could detect drug use.

Yet a big fraction lied about their drug use anyway.

Example: drug surveys

Drug surveys are really important:

- Drug abuse, whether it's cocaine in Detroit or bathtub speed in Nebraska, is a huge social problem.
- The problem fills our jails, drains public finances, and perpetuates a trans-generational cycle of poverty.
- Getting good data on this problem is important! Doctors, schools, and governments all rely on self-reported measures of drug use to guide their thinking on this issue.

Delaney-Black's asks: **can we trust any of it?**

It's not just drug surveys

Here are some other things that, according to research *on surveys*, people lie about *in surveys*.

Example: a better drug survey

Suppose that you want to learn about the prevalence of drug use among college students. Here's a cute trick that uses probability theory to mitigate someone's incentive to lie.

Suppose that, instead of asking people point-blank, you tell them:

- Flip a coin. Look at the result, but keep it private.
- If the coin comes up heads, please use the space provided to write an answer to question Q1: “Is the last digit of your tax ID number (e.g. SSN) odd?”
- If the coin comes up tails, please use the space provided to write an answer to question Q2: “Have you smoked marijuana in the last year?”

Example: a better drug survey

Key fact here: only the respondent knows which question he or she is answering.

- This gives people plausible deniability.
- Someone answering “yes”“ might have easily flipped heads and answered the first, innocuous question rather than the second, embarrassing one.
- The survey designer would never know the difference.

This reduces the incentive to lie.

Let's run the survey! Flip a coin, keep the result private, and then answer the question in public :-)

Analyzing the results

Notation:

- Let Y be the event “a randomly chosen subject answers yes.”
- Let Q_1 be the event “the subject answered question 1, about their tax ID number.”
- Let Q_2 be the event “the subject answered question 2, about marijuana use.”

By the rule of total probability:

$$\begin{aligned} P(Y) &= P(Y, Q_1) + P(Y, Q_2) \\ &= P(Q_1) \cdot P(Y | Q_1) + P(Q_2) \cdot P(Y | Q_2) \end{aligned}$$

Analyzing the results

$$P(Y) = P(Q_1) \cdot P(Y | Q_1) + P(Q_2) \cdot P(Y | Q_2)$$

$P(Y)$ is a weighted average of two conditional probabilities:

- $P(Y | Q_1)$, the probability that a subject answers “yes” when answering the tax-ID-number question.
- $P(Y | Q_2)$, the probability that a subject answers “yes” when answering the marijuana question.

This equation has five probabilities in it.

- **Which ones do we know from our survey?**
- **Which one do we care about?**

Analyzing the results

Let's solve for $P(Y | Q_2)$:

$$P(Y) = P(Q_1) \cdot P(Y | Q_1) + P(Q_2) \cdot P(Y | Q_2)$$

So

$$P(Y | Q_2) = \frac{P(Y) - P(Q_1) \cdot P(Y | Q_1)}{P(Q_2)}$$

Let's plug in our numbers on the right-hand side and get an answer! Feel free to chat with your neighbors.

Independence

Two events A and B are *independent* if

$$P(A \mid B) = P(A \mid \text{not } B) = P(A)$$

In words: A and B convey no information about each other:

- $P(\text{flip heads second time} \mid \text{flip heads first time}) = P(\text{flip heads second time})$
- $P(\text{stock market up} \mid \text{bird poops on your car}) = P(\text{stock market up})$
- $P(\text{God exists} \mid \text{Longhorns win title}) = P(\text{God exists})$

So if A and B are independent, then $P(A, B) = P(A) \cdot P(B)$.

Independence

Two events A and B are *conditionally independent*, given C , if

$$P(A, B \mid C) = P(A \mid C) \cdot P(B \mid C)$$

A and B convey no information about each other, once we know C :

$$P(A \mid B, C) = P(A \mid C).$$

Neither independence nor conditional independence implies the other.

- It is possible for two outcomes to be *dependent* and yet *conditionally independent*.
- Less intuitively, it is possible for two outcomes to be *independent* and yet *conditionally dependent*.

Independence

Let's see an example. Alice and Brianna live next door to each other and both commute to work on the same metro line.

- A = Alice is late for work.
- B = Brianna is late for work.

A and B are *dependent*: if Brianna is late for work, we might infer that the metro line was delayed or that their neighborhood had bad weather. This means Alice is more likely to be late for work:

$$P(A \mid B) > P(A)$$

Independence

Now let's add some additional information:

- A = Alice is late for work.
- B = Brianna is late for work.
- C = The metro is running on time and the weather is clear.

A and B are *conditionally independent*, given C. If Brianna is late for work but we know that the metro is running on time and the weather is clear, then we don't really learn anything about Alice's commute:

$$P(A \mid B, C) = P(A \mid C)$$

Independence

Same characters, different story:

- A = Alice has blue eyes.
- B = Brianna has blue eyes.

A and B are *independent*: Alice's eye color can't give us information about Brianna's.

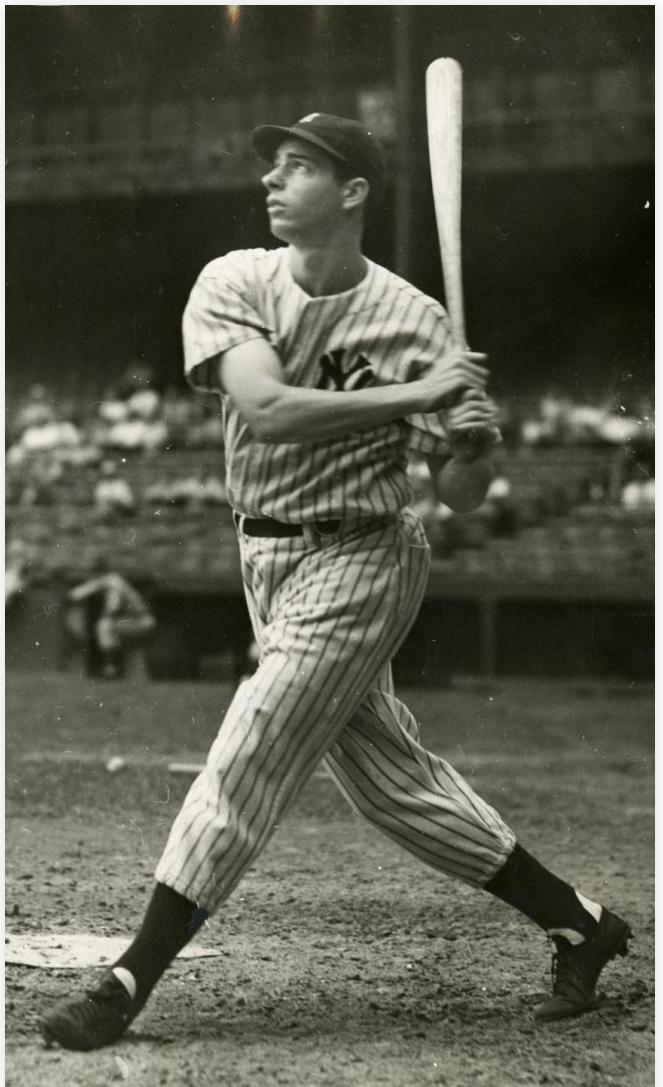
Independence

Again, let's add some additional information.

- A = Alice has blue eyes.
- B = Brianna has blue eyes.
- C = Alice and Brianna are sisters.

A and B are *conditionally dependent*, given C: if Alice has blue eyes, and we know that Brianna is her sister, then we know something about Brianna's genes. It is now more likely that Brianna has blue eyes.

Independence



Independence (or conditional independence) is often something we choose to assume for the purpose of making calculations easier. For example:

- Joe DiMaggio got a hit in about 80% of the baseball games he played in.
- Suppose that successive games are independent: if JD gets a hit today, it doesn't change the probability he's going to get a hit tomorrow.
- Then $P(\text{hit in game 1, hit in game 2}) = 0.8 \cdot 0.8 = 0.64$.

Independence

This works for more than two events. For example, Joe DiMaggio had a 56-game hitting streak in the 1941 baseball season. This was pretty unlikely!

$$\begin{aligned} & P(\text{hit game 1, hit game 2, hit game 3, ..., hit game 56}) \\ &= P(\text{hit game 1}) \cdot P(\text{hit game 2}) \cdot P(\text{hit game 3}) \cdots P(\text{hit game 56}) \\ &= 0.8 \cdot 0.8 \cdot 0.8 \cdots 0.8 \\ &= 0.8^{56} \\ &\approx \frac{1}{250,000} \end{aligned}$$

I like to call this the “compounding rule.”

Independence

Let's compare this with the corresponding probability for Pete Rose, a player who got a hit in 76% of his games. He's only *slightly* less skillful than DiMaggio! But:

$$\begin{aligned} & P(\text{hit game 1, hit game 2, hit game 3, \dots, hit game 56}) \\ &= 0.76^{56} \\ &\approx \frac{1}{5 \text{ million}} \end{aligned}$$

Small difference in one game, but a big difference over the long run.

Independence

What about an average MLB player who gets a hit in 68% of his games:

$$\begin{aligned} & P(\text{hit game 1, hit game 2, hit game 3, ..., hit game 56}) \\ &= 0.68^{56} \\ &\approx \frac{1}{2.5 \text{ billion}} \end{aligned}$$

Never gonna happen!

Independence

Summary:

- Joe DiMaggio: 80% one-game hit probability, 1 in 250,000 streak probability
- Pete Rose: 76% one-game hit probability, 1 in 5 million streak probability
- Average player: 68% one-game hit probability, 1 in 2.5 billion streak probability

A small difference in probabilities becomes an enormous gulf over the long term.

Lesson: probability compounds **multiplicatively**, like the interest on your credit cards.

What about everyday "hitting streaks"...?

- A mutual-fund manager outperforms the stock market for 15 years straight.
- A World-War II airman completes 25 combat missions without getting shot down, and gets to go home.
- A retired person successfully takes a shower for 1000 days in a row without slipping.
- A child goes 180 school days, or 1 year, without catching a cold from other kids at school. (Good luck!)

An example: surviving falls



Gerald Ford falls down the steps of Air Force One. (He survived.)

An example: surviving falls



Gerald Ford falls while skiing. (He again survived.)

An example: surviving falls



Gerald Ford falls at a summit in Salzburg. (He once again survived.)

An example: surviving falls

Chevy Chase playing Gerald Ford on Saturday Night Live

An example: surviving falls

- On an average day, 0.00003% of human beings will die in a fall.
- So the “average” daily fall-survival probability is 0.9999997.
- What about over 30 years? Assuming independence:

$$\begin{aligned}P(\text{30-year streak without a deadly fall}) &= (0.9999997)^{365 \times 30} \\&\approx 0.997\end{aligned}$$

So if you have an average daily risk, then you have a 0.3% chance of dying in a fall at some point over the next 30 years—hardly negligible, but still small.

An example: surviving falls

Let's change the numbers a tiny bit.

What if your daily survivorship probability was a bit smaller than average?

- Maybe you forgot to put a towel down on the bathroom floor after a shower?
- Maybe you never bother holding the handrail as you walk down the stairs?

To invoke the DiMaggio/Rose example: what if you became only slightly less skillful at not falling?

An example: surviving falls

For some specific numbers, let's make a diet analogy:

- Imagine that every day you have a single mid-morning Tic-Tac, which has 2 calories.
- One day, you decide to give it up.
- But you're wary of crash diets, so you decide to go slowly: you'll forego that Tic-Tac only once every 10 days.

An example: surviving falls

For some specific numbers, let's make a diet analogy:

- Imagine that every day you have a single mid-morning Tic-Tac, which has 2 calories.
- One day, you decide to give it up.
- But you're wary of crash diets, so you decide to go slowly: you'll forego that Tic-Tac only once every 10 days.

You've just reduced your average daily calorie consumption by about 1/100th of a percent. Will you lose weight over the long run? Probably not.

An example: surviving falls

But what if you reduced your daily fall-survivorship probability by 1/100 of a percent? (From 99.99997% to “merely” 99.99%).

Tiny change in the short run, big change in the long run:

$$\begin{aligned} P(\text{30-year streak without a deadly fall}) &= (0.9999)^{365 \times 30} \\ &\approx 0.33 \end{aligned}$$

Again: probability compounds multiplicatively (like interest), not additively (like calories).

Checking independence from data

Suppose we have two random outcomes A and B and we want to know if they're independent or not.

Solution:

- Check whether B happening seems to change the probability of A happening.
- That is, verify using data whether
$$P(A \mid B) = P(A \mid \text{not } B) = P(A)$$
- These probabilities won't be *exactly* alike because of statistical fluctuations, especially with small samples.
- But with enough data they should be pretty close if A and B are independent.

Example: the hot hand



NBA Jam c. 1993

Example: the hot hand

The “hot hand hypothesis” says that if a player makes their *previous* shot, they’re more likely to make their *next* shot (“He’s on fire!”):

$$P(\text{makes next} \mid \text{makes previous}) > P(\text{makes next} \mid \text{misses previous})$$

On the other hand, the “independence hypothesis” says that

$$P(\text{makes next} \mid \text{makes previous}) = P(\text{makes next} \mid \text{misses previous})$$

Example: the hot hand

The next slide show some data on shooting percentages for Dr. J's 1980–81 Philadelphia 76ers.

Key question: do players shoot better, worse, or about the same after they've just *made* a basket, versus how they do after they've just *missed* a basket?

Let's look at the data...



© NBAE/Getty Images

Example: the hot hand

Shooting percentages after:

| Player | 3 misses | 2 misses | 1 miss | overall | 1 hit | 2 hits | 3 hits |
|------------------|----------|----------|--------|---------|-------|--------|--------|
| Julius Erving | 0.52 | 0.51 | 0.51 | 0.52 | 0.52 | 0.53 | 0.48 |
| Caldwell Jones | 0.50 | 0.48 | 0.47 | 0.43 | 0.47 | 0.45 | 0.27 |
| Maurice Cheeks | 0.77 | 0.6 | 0.6 | 0.54 | 0.56 | 0.55 | 0.59 |
| Daryl Dawkins | 0.88 | 0.73 | 0.71 | 0.58 | 0.62 | 0.57 | 0.51 |
| Lionel Hollins | 0.50 | 0.49 | 0.46 | 0.46 | 0.46 | 0.46 | 0.32 |
| Bobby Jones | 0.61 | 0.58 | 0.58 | 0.47 | 0.54 | 0.53 | 0.53 |
| Andrew Toney | 0.52 | 0.53 | 0.51 | 0.40 | 0.46 | 0.43 | 0.34 |
| Clint Richardson | 0.50 | 0.47 | 0.56 | 0.50 | 0.50 | 0.49 | 0.48 |
| Steve Mix | 0.70 | 0.56 | 0.52 | 0.48 | 0.52 | 0.51 | 0.36 |

Which hypothesis looks right: hot hand or independence?
(Remember small-sample fluctuations.)

When independence goes wrong

Suppose we pick a random US family with four male children. What is the probability P that all four will be colorblind?

The probability that a randomly sampled US male is colorblind is about 8%. So the naive answer involves just compounding up this probability:

$$P = 0.08^4 \approx 0.0004$$

What's wrong here?



When independence goes wrong

Colorblindness runs in families (it's an X-linked trait, so males only need one copy on their X chromosome to express the phenotype). So it may be true that

$$P(\text{brother 1 colorblind}) = 0.08$$

But

$$P(\text{brother 2 colorblind} \mid \text{brother 1 colorblind}) = 0.5 \neq 0.08$$

And the same is true for all subsequent brothers: if brother 1 is colorblind, you know that mom is a carrier, and so all her male children have a 50/50 chance of colorblindness (conditional independence, given mom's genes!)

When independence goes wrong

The correct overall probability has to be built up piece by piece using the multiplication rule:

$$\begin{aligned} P(\text{brothers 1-4 colorblind}) &= P(\text{brother 1 colorblind}) \\ &\quad \times P(\text{brother 2 colorblind} \mid \text{brother 1 colorblind}) \\ &\quad \times P(\text{brother 3 colorblind} \mid \text{brothers 1-2 colorblind}) \\ &\quad \times P(\text{brother 4 colorblind} \mid \text{brothers 1-3 colorblind}) \end{aligned}$$

So:

$$P(\text{brothers 1-4 colorblind}) = 0.08 \times 0.5^3 = 0.01$$

When independence goes wrong

Seems silly, right?

But you'd be surprised at how often people make this mistake! We might call this the “fallacy of mistaken compounding”: assuming events are independent and naively multiplying their probabilities.

Out of class, I'm asking you to read two short pieces that illustrate this unfortunate reality:

- How likely is it that birth control could let you down? from the *New York Times*
- An excerpt from Chapter 7 of [AIQ: How People and Machines are Smarter Together](#), by Nick Polson and James Scott.

Bayes' Rule

Key fact: all probabilities are contingent on what we know.

When our knowledge changes, our probabilities must change, too.

Bayes' rule tells us how to change them. Suppose A is some event we're interested in and B is some new relevant information. Bayes' rule tells us how to move from a prior probability, $P(A)$, to a posterior probability $P(A | B)$ that incorporates our knowledge of B.

Bayes' Rule

$$P(A \mid B) = P(A) \cdot \frac{P(B \mid A)}{P(B)}$$

- $P(A)$ is the prior probability: how probable is A , before having seen data B ?
- $P(A \mid B)$ is the posterior probability: how probable is A , now that we've seen data B ?
- $P(B \mid A)$ is the likelihood: if A were true, how likely is it that we'd see data B ?
- $P(B)$ is the marginal probability of B : how likely is it that we'd see data B overall, regardless of whether A is true or not?

Calculating $P(B)$: use the rule of total probability.

Bayes' Rule: a toy example

Imagine a jar with 1024 normal quarters. Into this jar, a friend places a single two-headed quarter (i.e. with heads on both sides). Your friend shakes the jar to mix up the coins. You draw a single coin at random from the jar, and without examining it closely, flip the coin ten times.

The coin comes up heads all ten times.

Are you holding the two-headed quarter, or an ordinary quarter?

Not just a toy example: in any industry where companies compete strenuously for talent, a lot of time and energy is spent looking for “two-headed quarters”!

Bayes' Rule: a toy example

A real-world version of the two-headed quarter problem:

- Suppose you're in charge of a large trading desk at a major Wall Street bank.
- You have 1025 employees under you, and each one is responsible for managing a portfolio of stocks to make money for your firm and its clients.
- One day, a young trader knocks on your door and confidently asks for a big raise. You ask her to make a case for why she deserves one.

Bayes' Rule: a toy example

She replies:

Look at my record: I'm the best trader on your floor. I've been with the company for ten months, and in each of those ten months, my portfolio returns have been in the top half of all the portfolios managed by my peers on the trading floor. If I were just one of those other average Joes, this would be very unlikely. In fact, the probability that an average trader would see above-average results for ten months in a row is only $(1/2)^{10}$, which is less than one chance in a thousand. Since it's unlikely I would be that lucky, I should get a raise.

Bayes' Rule: a toy example

Is the trader lucky, or good? Same math as the big jar of quarters!

- Metaphorically, the trader is claiming to be the two-headed coin (T) in a sea of mediocrity.
- Her data is “D = ten heads in a row”: she's performed above average for ten months straight.
- This is admittedly unlikely: $P(D | T) = 1/2^{10} = 1/1024$.
- But excellent performers are probably also rare, so that the prior probability $P(T)$ is pretty small to begin with.

To make an informed decision, you need to know $P(T | D)$: the posterior probability that the trader is an above-average performer, given the data.

Bayes' Rule: a toy example

So let's return to the two-headed quarter example and see how a posterior probability is calculated using Bayes' rule:

$$P(T \mid D) = \frac{P(T) \cdot P(D \mid T)}{P(D)}.$$

We'll take this equation one piece at a time.

Bayes' Rule: a toy example

$$P(T \mid D) = \frac{P(T) \cdot P(D \mid T)}{P(D)}.$$

$P(T)$ is the prior probability that you are holding the two-headed quarter.

- There are 1025 quarters in the jar: 1024 ordinary ones, and one two-headed quarter.
- Assuming that your friend mixed the coins in the jar well enough, then you are just as likely to draw one coin as another.
- So $P(T)$ must be 1/1025.

Bayes' Rule: a toy example

$$P(T | D) = \frac{P(T) \cdot P(D | T)}{P(D)}.$$

Next, what about $P(D | T)$, the likelihood of flipping ten heads in a row, given that you chose the two-headed quarter?

- Clearly this is 1.
- If the quarter has two heads, there is no possibility of seeing anything else.

Bayes' Rule: a toy example

$$P(T \mid D) = \frac{P(T) \cdot P(D \mid T)}{P(D)}.$$

Finally, what about $P(D)$, the marginal probability of flipping ten heads in a row? Use the rule of total probability:

$$P(D) = P(T) \cdot P(D \mid T) + P(\text{not } T) \cdot P(D \mid \text{not } T).$$

- $P(T)$ is $1/1025$, so $P(\text{not } T)$ is $1024/1025$.
- $P(D \mid T) = 1$.
- $P(D \mid \text{not } T)$ is the probability of a ten-heads “winning streak”:

$$P(D \mid \text{not } T) = \left(\frac{1}{2}\right)^{10} = \frac{1}{1024}.$$

Bayes' Rule: a toy example

We can now put all these pieces together:

$$\begin{aligned} P(T \mid D) &= \frac{P(T) \cdot P(D \mid T)}{P(T) \cdot P(D \mid T) + P(\text{not } T) \cdot P(D \mid \text{not } T)} \\ &= \frac{\frac{1}{1025} \cdot 1}{\frac{1}{1025} \cdot 1 + \frac{1024}{1025} \cdot \frac{1}{1024}} = \frac{1/1025}{2/1025} \\ &= \frac{1}{2}. \end{aligned}$$

There is only a 50% chance that you are holding the two-headed coin. Yes, flipping ten heads in a row with a normal coin is very unlikely (low likelihood). But so is drawing the one two-headed coin from a jar of 1024 normal coins! (Low prior probability.)